# Diffuman4D: 4D Consistent Human View Synthesis from Sparse-View Videos with Spatio-Temporal Diffusion Models

Yudong Jin[1]    Sida Peng[1]    Xuan Wang[2]    Tao Xie[1]    Zhen Xu[1]
Yifan Yang[1]    Yujun Shen[2]    Hujun Bao[1]    Xiaowei Zhou[1†]
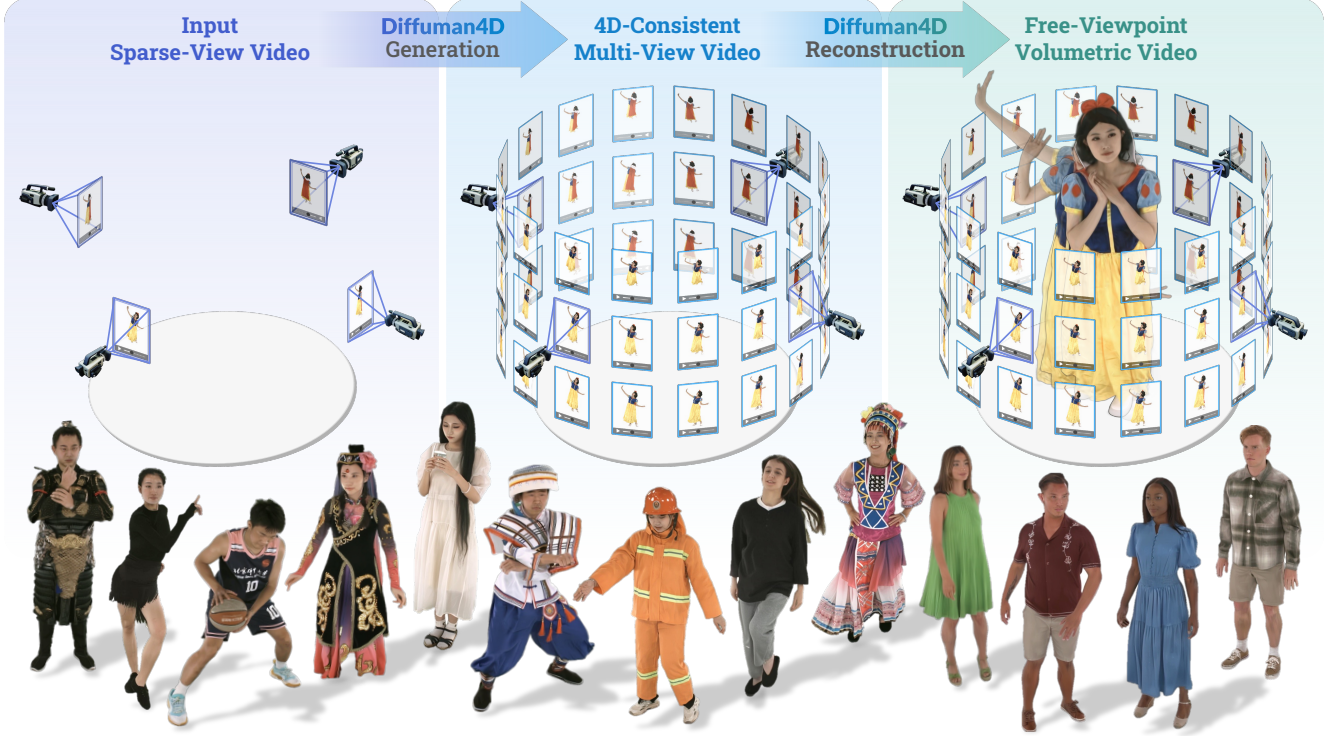
[1]Zhejiang University    [2]Ant Research

Figure 1. Diffuman4D enables high-fidelity free-viewpoint rendering of human performances from sparse-view videos. The bottom row shows representative results.

## Abstract

*This paper addresses the challenge of high-fidelity view synthesis of humans with sparse-view videos as input. Previous methods solve the issue of insufficient observation by leveraging 4D diffusion models to generate videos at novel viewpoints. However, the generated videos from these models often lack spatio-temporal consistency, thus degrading view synthesis quality. In this paper, we propose a novel sliding iterative denoising process to enhance the spatio-temporal consistency of the 4D diffusion model. Specifically, we define a latent grid in which each latent encodes the image, camera pose, and human pose for a certain viewpoint and timestamp, then alternately denoising the latent grid along spatial and temporal dimensions with a sliding window, and finally decode the videos at target viewpoints from the corresponding denoised latents. Through the iterative sliding, information flows sufficiently across the latent grid, allowing the diffusion model to obtain a large receptive field and thus enhance the 4D consistency of the output, while making the GPU memory consumption affordable. The experiments on the DNA-Rendering and ActorsHQ datasets demonstrate that our method is able to synthesize high-quality and consistent novel-view videos and significantly outperforms the existing approaches. See our project page for interactive demos and video results: https://diffuman4d.github.io/.*

† Corresponding author: Xiaowei Zhou.

# 1. Introduction

This paper aims to address the problem of high-quality 4D view synthesis of humans in motion from sparse-view videos, which has wide applications in augmented reality, film production, sports broadcasting, etc. Traditional multi-view stereo-based methods [17, 53, 54] and recent neural rendering-based methods [67, 72, 76] generally require a dense array of synchronized cameras to capture a sufficient number of views for high-quality reconstruction, making them difficult to apply to real-world scenarios. When input views become sparse, these methods tend to fail as the insufficient observations make the reconstruction problem ill-posed.

An intuitive solution to this problem is to leverage conditional image or video generative models to generate novel views conditioned on the input views [65, 68, 69]. By leveraging the attention mechanism, these methods inject spatial and temporal control signals into video generative models, aiming to produce human images at target viewpoints and timestamps. However, these methods often struggle with the spatio-temporal consistency of generated images, especially when the human topology and cloth deformations are complex. A key reason is that, due to GPU memory limitations, these methods often have to generate the target images in multiple passes, and the inherent probabilistic nature of generative models leads to variance among the outputs.

In this paper, we propose a novel spatio-temporal diffusion model for generating 4D consistent multi-view human videos. Our key innovation lies in a novel sliding iterative denoising process that ensures the 4D consistency of our model's outputs. Specifically, given a sparse-view video, we first encode image observations and camera parameters as conditioning latents, and define noise latents at target viewpoints and timestamps, forming a 4D latent grid. To denoise the latent grid, we use a window that alternately slides forward and backward in the spatial and temporal dimensions. In contrast to the previous method [68] that executes the full denoising process at each sliding operation, our model only performs a few denoising steps during sliding. This iterative sliding strategy ensures sufficient information propagation across the latent grid, enabling the diffusion model to incorporate surrounding 4D signals to produce each target output, while dynamically adjusting their influence based on spatio-temporal distance.

To further boost the consistency of our generative model, we exploit the 3D human skeleton sequence as a structural prior to guide the generation process. Concretely, our method first extracts a 3D human skeleton sequence from the given sparse-view video using an off-the-shelf 3D pose estimator. Then, we project the skeletons into the image space at each viewpoint and each timestamp, serving as the conditioning signal alongside camera parameters and image observations to guide the spatio-temporal diffusion model.

Finally, we decode the videos at target viewpoints from the corresponding denoised latents, and reconstruct high-quality 4D Gaussian Splatting (4DGS) [67, 73, 76] based on the input-view and the synthesized novel-view videos.

To enable model training, we meticulously process the DNA-Rendering [10] dataset by recalibrating camera parameters, optimizing image color correction matrices (CCMs), predicting foreground masks, and estimating human skeletons. We compare our method against state-of-the-art methods on DNA-Rendering and ActorsHQ datasets, highlighting the performance of our framework in capturing detailed human motion and appearance from sparse-view inputs.

In summary, our contributions are as follows:
- We introduce Diffuman4D, a novel diffusion model that generates spatio-temporally consistent and high-resolution (1024p) human videos from sparse-view video inputs.
- We propose a sliding iterative denoising mechanism that enhances both spatial and temporal consistency of generated long-term videos while maintaining efficient inference.
- We design a human pose conditioning scheme to enhance appearance quality and motion accuracy of generated human videos.
- We plan to release our processed version of the DNA-Rendering dataset, which we believe will benefit future research in this area.

# 2. Related Work

**4D reconstruction from dense views.** Reconstructing dynamic 3D human performances from captured videos and performing novel view synthesis to create immersive playbacks has been a long-standing research problem in computer vision and graphics. Traditional methods tend to leverage complicated hardware, such as dense camera arrays [11, 19, 25, 59, 60] or depth sensors [2, 5, 43, 56, 61], to reconstruct high-fidelity human performances. Recent advancements in neural scene representation, namely neural radiance field (NeRF) [42] and 3D Gaussian Splatting (3DGS) [26], have demonstrated remarkable success in static 3D scene reconstruction. [7, 14, 16, 32, 55, 63, 72, 73, 76] propose to lift this base 3D representation (NeRF or 3DGS) to 4D by adding an additional temporal dimension, enabling it to model temporal variations in dynamic scenes. However, these methods still heavily rely on densely distributed and well-synchronized multi-view video inputs and exhibit severe overfitting under sparse-view conditions, which greatly limits their accessibility.

**4D reconstruction from sparse views.** To alleviate the requirement of dense multi-view inputs, some methods [48, 66] leverage human priors such as SMPL [40] to

guide the reconstruction process. [22, 47, 71] propose constructing a 3D static canonical space with a neural field and learning a deformation field [15, 33, 45, 46, 50] based on the SMPL human prior to map dynamic elements back to this canonical space. While effective in certain scenarios, these approaches face challenges in accurately estimating shape deformations, particularly when dealing with complex outfits or rapid motions. On the other hand, methods like [8, 9, 34, 38, 85] explore the use of depth priors, such as stereo [36] or multi-view [77, 78] depth estimation for generalizable scene reconstruction and novel view synthesis. However, these methods are highly dependent on the accuracy of depth estimation, often struggling in cases involving occlusions, textureless regions, or extremely sparse viewpoints.

**4D generation.** Recent developments in 3D content generation [18, 20, 37, 39, 49, 57, 64] and video diffusion technologies [6, 21, 29, 75] provide a promising direction for handling challenging scenarios mentioned above by introducing generative data priors to the reconstruction pipeline. [3, 4, 51, 58, 79, 81, 83, 86] leverage the Score Distillation Sampling (SDS) [49] to extract 4D representations from image or video diffusion models. However, SDS limits their scalability to large-scale 4D reconstruction tasks due to its high computational cost and tends to produce oversmoothed geometries and unrealistic textures. To overcome these limitations, [31, 44, 69, 74, 80] propose to condition the diffusion models to generate spatio-temporal consistent multi-view videos, which can be further used for 4D reconstruction. However, these methods only focus on object-level generation. Recently, CAT4D [68] has made a step forward by proposing a multi-view consistent video diffusion model that can handle general scenes by leveraging general time embedding and Plücker embedding [24] for temporal and spatial consistency. Despite the improvement, CAT4D still struggles with complex human generation due to the heavy shape distortion and self-occlusions caused by the motion of soft structures like hair and clothes, and it is difficult for diffusion models to solve these ill-posed problems by only relying on the general conditions. In this paper, we propose to introduce additional human-specific priors to address these challenges.

## 3. Method

We reconstruct human performances from sparse-view videos in two stages. First, we transform the input sparse-view videos into dense multi-view videos using our spatio-temporal diffusion model. Then, we reconstruct the human performance by optimizing a 4D Gaussian Splatting (4DGS) from these generated multi-view videos. We first describe our spatio-temporal diffusion model (Sec. 3.1) and the denoising mechanism for generating spatio-temporal

consistent multi-view videos (Sec. 3.2). Then, we describe the skeleton conditioning scheme (Sec. 3.3), and the method we used to reconstruct human performances (Sec. 3.4).

### 3.1. Spatio-Temporal Diffusion Model

**Pipeline.** As illustrated in Fig. 2, our spatio-temporal diffusion model takes $M$-view videos as input and aims to generate $N$-view target videos, where all videos consist of $T$ frames. We first encode the input videos into a latent space using a pretrained VAE and generate noise latents for the target videos. These latents are structured into a sample grid of size $(N + M) \times T$, where the two axes represent the spatial (multi-view) dimension and the temporal (video) dimension. Each sample within the grid comprises an input image latent (or a target noise latent) along with its corresponding conditioning embeddings, which consist of a skeleton latent and Plücker coordinates (see Sec. 3.3 for details). We then employ a sliding iterative approach to progressively denoise the sample grid (see Sec. 3.2 for details). Next, we decode the target image latents into corresponding videos. Finally, we reconstruct a high-quality 4D Gaussian Splatting (4DGS) representation of the human performance using both the input-view and target-view videos, enabling real-time rendering.

**Architecture.** Our model follows the architecture of multi-view latent diffusion models [18, 57]. Specifically, the model employs 3D self-attention layers to enable information exchange across images. The images are encoded into latent representations through a pretrained VAE, allowing the model to learn the joint distribution in the latent space. We disable the text conditioning by setting the input prompt to an empty string.

### 3.2. Sliding Iterative Denoising

To achieve high-quality 4D human reconstruction, the input data must be sufficiently dense in both spatial (multi-view) and temporal (video) dimensions. For instance, reconstructing a 10-second 4D human typically requires tens of thousands of input images. Due to GPU memory constraints, existing video diffusion models can only denoise a limited number of images per inference, requiring the entire image sequence to be split into hundreds of groups. Since these denoising iterations operate independently, the generated images exhibit inconsistency due to variations in the diffusion denoising process. The recent approach [68] introduces a sliding-window strategy that separately complete the denoising process within each sliding window and then apply median filtering to the overlapping samples, thereby reducing the variance of diffusion denoising. However, this method still suffers from inconsistency when generating long-sequence samples, and performing multiple denoising iterations significantly increases the inference time. See Fig. 5 for comparative results.
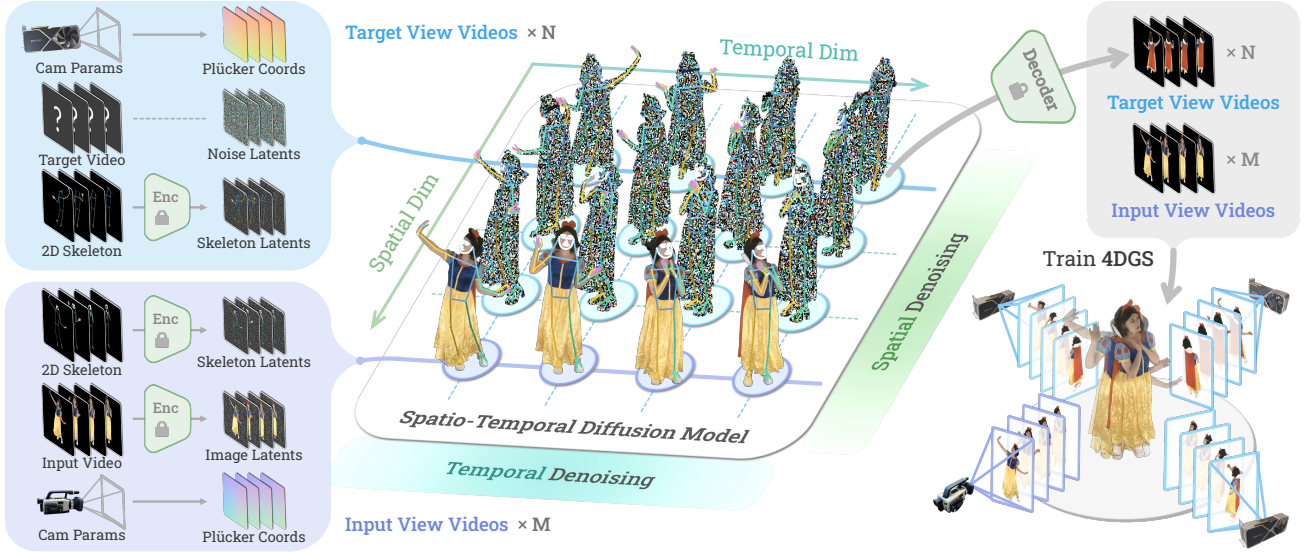
Figure 2. **Overview of Diffuman4D**. Our model takes $M$ input-view videos as input, generates $N$ target-view videos, and reconstructs a 4DGS of the human using both input and generated videos. Specifically, the input-view videos are encoded into the latent space using a pretrained VAE. The 3D human-skeleton sequence is projected into each view, rendered as RGB maps, and encoded into the same latent space. In addition, the camera parameters are encoded as Plücker coordinates [82]. The skeleton latents and Plücker coordinates are then concatenated with the image latents at input views or the noise latents at target views, forming input samples or target samples, respectively. The samples across all views and timestamps form a sample grid and are fed into our spatio-temporal diffusion model that denoises the samples using a sliding iterative denoising mechanism (see Sec. 3.2). Then, the denoised image latents for the target views are decoded into target-view videos using a pretrained VAE decoder. Finally, the 4DGS is reconstructed using an off-the-shelf approach [73] (see Sec. 3.4), enabling real-time novel view rendering.

**Sliding iterative denoising.** To address the aforementioned challenges, we propose a sliding iterative denoising mechanism that enhances consistency in long-sequence generation by leveraging rich contextual information during the denoising process. Specifically, given a target sample sequence, we set a context window of length $W$ that slides over the sequence with a fixed stride $S$. In each iteration, we concatenate the target samples with input samples and feed them into the diffusion model for denoising $P$ steps.

As illustrated in Fig. 3a, given a human-centric sample sequence (e.g., camera views arranged circularly), we first slide the context window counter-clockwise to enable information propagation along that direction. Subsequently, we reverse the sliding direction to clockwise, allowing bidirectional context aggregation. Formally, each sample is denoised $D = 2 \times P \times W/S$ steps in total. We accordingly set the diffusion inference steps to $D$, therefore the generation is completed after the above process.

The same operation can be applied to the temporal sample sequence, enabling each sample to aggregate both past and future contexts during denoising. It is worth noting that multiple sample sequences can be denoised in parallel using multiple GPUs. The sliding iterative denoising mechanism allow the diffusion model to harmonize the consistency of the target samples in both spatial and temporal dimensions, enabling the generation of high-quality 4D image grids.

**Alternating denoising.** Following the previous work [68], we adopt an alternating denoising strategy to further improve the spatial and temporal consistency of generated images. Fig. 3b illustrates our denoising process. Given an $M$-view, $T$-frame video, we first denoise target samples in the spatial dimension for $D/2$ steps, conditioning on the $M$ input views at the corresponding time, and then perform denoising in the temporal dimension for the remaining $D/2$ steps, conditioning on the $W$ frames within the corresponding time range from the nearest-distance view. This ensures the generation of spatio-temporally consistent samples.

By combining this strategy with sliding iterative denoising, each sample perceives surrounding spatial information from adjacent camera views and temporal information from neighboring video frames through sliding context windows. Furthermore, each target sample acts as a local center, and samples closer to it undergo more joint denoising steps within the context window. This mechanism aligns with the nature of 4D data: closer samples exhibit stronger correlations, which require more intensive consistency constraints.

### 3.3. Skeleton-Conditioned Diffusion

Previous 4D generation methods [65, 68] were designed for general scenes by directly injecting spatial signals (camera embeddings) and temporal signals (timestamp embeddings) into diffusion models. However, these methods face signif-
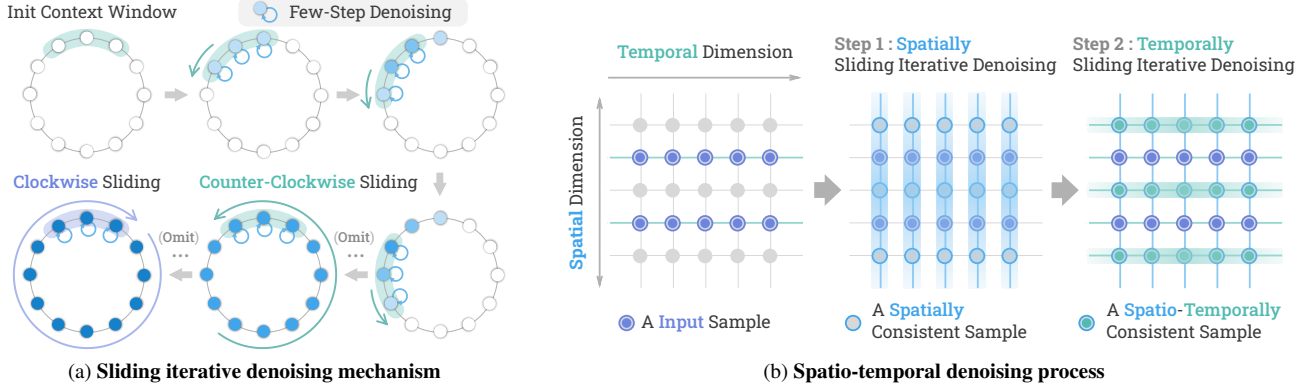
|                          |                          |
|:------------------------:|:------------------------:|
| (a) **Sliding iterative denoising mechanism** | (b) **Spatio-temporal denoising process** |

Figure 3. **Illustration of our denoising framework.** (a) **Sliding iterative denoising mechanism**. Given a view sequence arranged in a circular manner, we initialize a context window of length $W$ (here, $W = 3$). The window slides counterclockwise with a stride of $S$ (here, $S = 1$), and in each iteration, the samples within the context window are fed into the model for $P$ denoising steps. After completing a full rotation, we reverse the direction of the window, allowing it to slide clockwise for another full cycle. Consequently, each sample undergoes a total of $D = 2PW/S$ denoising steps. (b) **Spatio-temporal denoising process**. Given an $M$-view, $N$-frame video (here, $M = 2$, $N = 5$) and diffusion inference steps $D$, the model first performs denoising in the spatial dimension for $D/2$ steps, producing spatially consistent samples (latents). It then performs denoising in the temporal dimension for the remaining $D/2$ steps, resulting in spatio-temporally consistent samples. Sliding iterative denoising is employed in both spatial and temporal denoising processes. Each column or row can be denoised in parallel.

icant challenges when generating spatio-temporally consistent human images. First, while Plücker coordinates provide pixel-wise camera pose signals, the generated images often exhibit noticeable pose errors from the input cameras. Second, human motion often involves substantial deformation (such as hair and loose clothing), leading to various shape distortions and self-occlusions. Diffusion models struggle to stabilize the generation process.

**Skeleton conditioning.** To address these challenges, we introduce additional human-specific conditioning signals to constrain the model's generation space, leading to more accurate and consistent human image synthesis. Considering the vast diversity of human data (such as different clothing, body shapes, and genders), we require an intermediate human representation that ensures precise conditioning signals for the model. 3D human skeleton sequence emerges as an ideal choice to meet these requirements, as it can be easily recovered from sparse-view videos and provides a consistent 4D human structure. These 3D skeletons can be projected into 2D space, capturing a snapshot from a specific camera at a given timestamp, to enhance both temporal and spatial conditioning signals for diffusion models.

Specifically, we first estimate 2D human skeletons using Sapiens [27] and triangulate them to obtain a 3D skeleton sequence. We then project these skeletons into each view and render them as an RGB map, assigning different colors to different body parts to enrich the conditioning information. This RGB map is then encoded into the latent space using a pretrained VAE, serving as a pixel-aligned feature that significantly enhances generation quality for complex human poses.

**Skeleton-Plücker mixed conditioning.** However, skeleton predictions are often incomplete for individuals wearing complex clothing, leading to a degradation in the completeness of pose control signals. Additionally, since skeleton representations do not contain explicit occlusion information, they introduce inherent ambiguities of front-back symmetry. See Fig. 4 for an example. To mitigate these limitations, we retain Plücker coordinate conditioning to provide explicit camera pose information, thereby improving the robustness of the generation process.

### 3.4. 4DGS Reconstruction

Based on the approach above, our model can generate dense-view videos with spatio-temporal consistency. The output videos can be fed into any existing 4D reconstruction pipeline to yield a 4D human representation. In our experiments, we employ LongVolcap [70, 73] as our reconstruction method, which is an enhanced version of 4DGS that can effectively reconstruct a long volumetric video with a temporally hierarchical Gaussian representation.

## 4. Experiments

### 4.1. Implementation Details

We train the model on either spatial sample sequences or temporal sample sequences, each with a total length of $M + N = 16$. For the spatial sample sequence, all samples are captured simultaneously from different cameras, and the model is trained to generate $N = 12$ target samples from $M = 4$ conditional samples. For the temporal sample sequence, the model is trained to generate $N = 8$ consecutive

samples from a target camera $C_T$, conditioned on $M = 8$ samples over the same time range from a reference camera $C_R$, which is randomly selected during training. We train the model on each of the above sequences with an equal probability of 50%. Both spatial and temporal training sequences are randomly sampled from the entire spatial (temporal) candidates. During training, we randomly drop all conditions (including image latents, skeleton latents, and Plücker coordinates) with a probability of 10% to enable classifier-free guidance. Please refer to the supplementary materials for more details.

## 4.2. Datasets and Baselines

**Datasets.** We train our model on the DNA-Rendering [10] dataset, which contains over $2,000$ sequences of human performances in diverse outfits and dynamic motions. For training, we first filter out actors interacting with complex objects and then select $1,000$ sequences, each with 48 views and 225 frames per view, totaling 10 million images. We use 16 sequences from the test set covering a variety of clothing types and action categories for quantitative comparison. Additionally, we evaluate our model on the ActorsHQ [23] dataset, which consists of 12 sequences of human performances, to assess zero-shot generalization.

**Baselines.** We compare our approach with multiple categories of state-of-the-art methods: the optimization-based method LongVolcap [73], the SMPL-based method GauHuman [22], the feed-forward method GPS-Gaussian [85], and the generation-based method CAT4D† [68]. † indicates our reproduced version.

Since neither the DNA-Rendering nor ActorsHQ provides SMPL models for our test sequences, we use Easy-Mocap [1] to extract SMPL models for test sequences, which serve as the inputs for GauHuman [22]. For GPS-Gaussian [85], we follow their view-selection strategy: select the two input views closest to the target view as input to generate each view. We trained CAT4D† [68] on the processed DNA-Rendering dataset under the same training settings as our approach. CAT4D† [68] chooses conditional views from sparse-view videos to denoise each row or column, using the same sampling sequences and conditional-view selection strategy described in Sec. 3.2.

## 4.3. Comparison to Baselines

**Comparison on the DNA-Rendering [10] dataset.** We provide quantitative and qualitative comparisons on the DNA-Rendering [10] dataset in Tab. 1 and Fig. 6 respectively. As shown in the visualization and metrics, our method consistently outperforms baselines with higher visual quality and better spatio-temporal consistency. The optimization-based method LongVolcap [73] struggles with the ill-posed nature of sparse-view reconstruction, result-
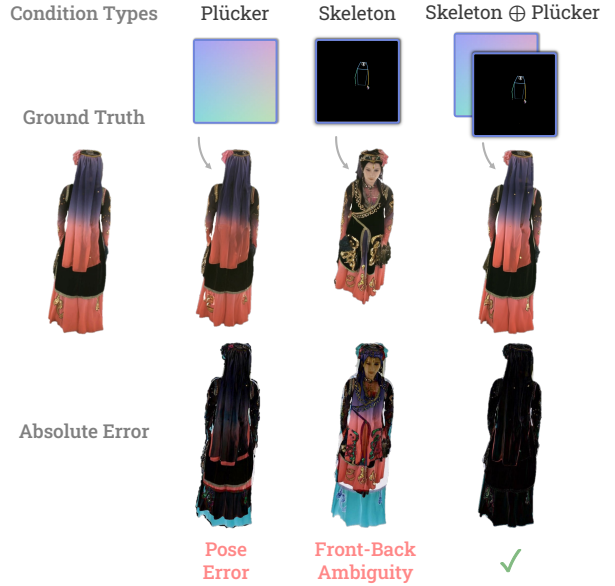


Figure 4. **Qualitative comparison of different conditioning**. The skeleton-Plücker mixed conditioning serves as a robust human pose prior for diffusion models.

ing in noisy renderings. GauHuman [22] fails to reliably reconstruct performers wearing complex clothing and executing dynamic motions. The depth estimator of GPS-Gaussian [85] breaks down under our sparse-view settings, yielding fragmented results on highly dynamic sequences. In contrast, our method not only effectively handles the challenging sparse-view setting by producing reasonable guidance from the diffusion prior, but can also generalize well to the complex motions and appearances of dynamic human performers. Note that even with only 4 input views, our method achieves visual quality comparable to the dense reconstruction from 48 views using LongVolcap [73].

**Comparison on the ActorsHQ [23] dataset.** As shown in Tab. 1 and Fig. 6, our model generalizes well to the unseen actor appearances and motions of the ActorsHQ dataset. In comparison, the baseline methods struggle to produce coherent geometry and appearance, consistent with their results on the DNA-Rendering dataset. Thanks to our unique model design, even with limited observations on the ActorsHQ dataset, Diffuman4D consistently produces much sharper and also spatio-temporally consistent human performance reconstruction results.

## 4.4. Ablation Study

We perform two ablation studies on the DNA-Rendering dataset [10], using three challenging motion sequences to evaluate the sliding iterative denoising mechanism and six complex-clothing sequences to assess the skeleton–Plücker conditioning scheme.

Table 1. **Quantitative comparison on DNA-Rednering [10] and ActorsHQ [23].** Diffuman4D surpasses baseline approaches across different settings and metrics. Note that CAT4D[†] is our reproduced version.

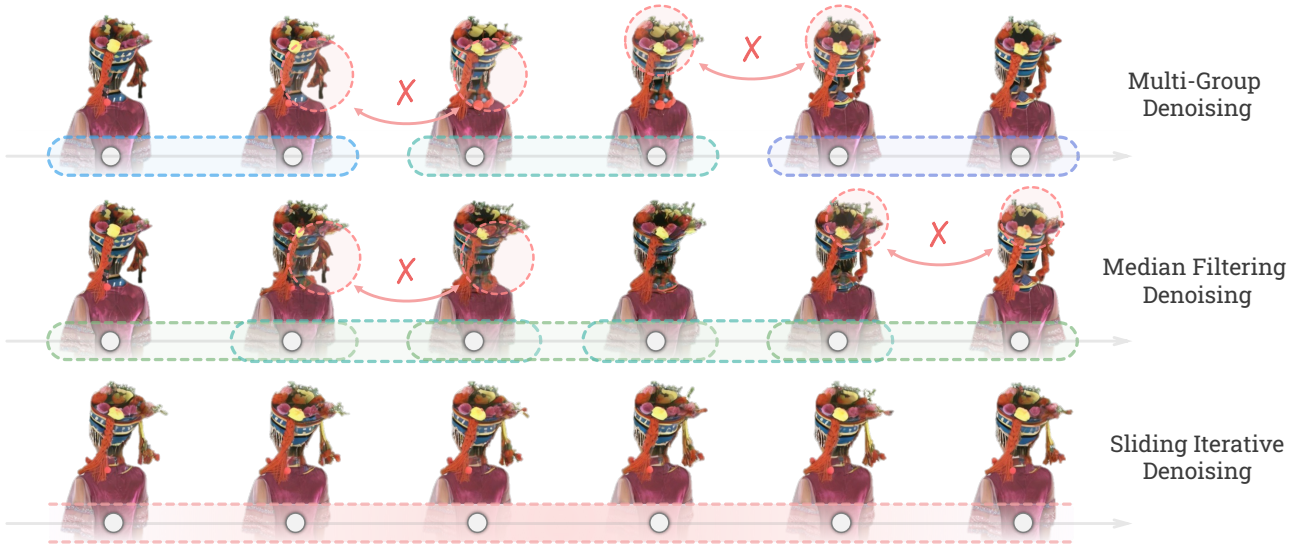| Dataset | | 4-View Video | | | 8-View Video | | |
|---|---|---|---|---|---|---|---|
| Method | Type | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ |
| **DNA-Rendering [10]** | | | | | | | |
| LongVolcap (4DGS) [73] | Optimization | 20.064 | 0.740 | 0.296 | 24.211 | 0.840 | 0.221 |
| GauHuman [22] | SMPL-Prior | 18.406 | 0.723 | 0.327 | 18.818 | 0.737 | 0.316 |
| GPS-Gaussian [85] | Feed-Forward | 11.250 | 0.457 | 0.460 | 17.604 | 0.714 | 0.270 |
| CAT4D[†] [68] | Generation | 21.445 | 0.806 | 0.234 | 22.531 | 0.824 | 0.221 |
| Ours | Generation | **25.393** | **0.864** | **0.161** | **26.324** | **0.881** | **0.150** |
| **ActorsHQ [23]** | | | | | | | |
| LongVolcap (4DGS) [73] | Optimization | 21.313 | 0.761 | 0.271 | 28.120 | 0.896 | 0.156 |
| GauHuman [22] | SMPL-Prior | 20.449 | 0.776 | 0.275 | 21.454 | 0.803 | 0.252 |
| GPS-Gaussian [85] | Feed-Forward | 10.562 | 0.453 | 0.481 | 14.208 | 0.601 | 0.379 |
| CAT4D[†] [68] | Generation | 21.562 | 0.808 | 0.229 | 23.002 | 0.873 | 0.206 |
| Ours | Generation | **27.875** | **0.903** | **0.121** | **28.747** | **0.916** | **0.110** |



Figure 5. **Qualitative comparison of different denoising strategies**. Our sliding iterative denoising method ensures a consistent appearance throughout a long image sequence.

Table 2. **Quantitative comparison of denoising strategies.**

| Sampling Method | PSNR↑ | SSIM↑ | LPIPS↓ |
|---|---|---|---|
| Multi-group | 20.913 | 0.753 | 0.224 |
| Median filtering | 21.609 | 0.766 | 0.226 |
| Sliding iterative | **22.363** | **0.778** | **0.196** |

Table 3. **Quantitative comparison of conditioning schemes.**

| Condition | PSNR↑ | SSIM↑ | LPIPS↓ |
|---|---|---|---|
| w/o skeleton | 16.864 | 0.608 | 0.272 |
| w/o Plücker | 20.753 | 0.638 | 0.203 |
| Ours | **22.120** | **0.707** | **0.184** |

**Diffusion denoising strategy.** We compare three diffusion sampling strategies: the multi-group denoising approach that divides the data into multiple groups, the median filtering approach that takes the median values of the denoised results of overlapping images between groups, and our sliding iterative diffusion denoising strategy. As shown in Fig. 5 and Tab. 2, the multi-group denoising approach does not take the temporal and spatial correlation between groups into consideration, leading to sudden jumps in the generated results between different segments. The median filtering approach slightly mitigates the intra-group inconsistency of the multi-group approach by taking the median values of the denoised results between different windows. However, the computational cost of this approach is inversely proportional to the overlap ratio, and could still produce incon-

| LongVolcap | GauHuman | GPS-Gaussian | CAT4D† | Ours | Ground Truth | LongVolcap | GauHuman | GPS-Gaussian | CAT4D† | Ours | Ground Truth |

(a) **Results on DNA-Rendering [10] test set.**  (b) **Zero-shot generalization on ActorsHQ [23].**

Figure 6. **Qualitative comparison on DNA-Rendering [10] and ActorsHQ [23].** GPS-Gaussian uses 8 input views while all other methods use 4 input views. CAT4D† is our reproduced version. Our Diffuman4D consistently outperforms baselines with higher visual quality and better spatio-temporal consistency.

sistency if the window doesn't overlap enough. In contrast, our sliding iterative denoising strategy introduces a smoothness-inductive bias during the denoising process of the diffusion model, at the same time maintaining constant computational cost by merging the sliding operation with the denoising steps. This process produces more consistent and globally accurate results compared to the two solutions.

**Conditioning scheme.** As shown in Fig. 4 and Tab. 3, we compare three conditioning schemes: w/o skeleton, w/o Plücker, and our skeleton-Plücker. The w/o-skeleton conditioning has limited camera control over the generated content, producing large misalignment due to the ill-posed nature of the generative problem. The w/o-Plücker conditioning can provide fine-grain control over the generated human, but it may struggle to disambiguate between the front and back, left and right of the generation target, producing inconsistent guidance for the reconstruction module. In comparison, our conditioning scheme combines the merits of the camera and pose control signal of the Plücker and skeleton embedding, generating consistent and controllable novel view results of the target human actor.

## 5. Conclusion

This paper introduces Diffuman4D, a novel diffusion model capable of generating high-resolution (1024p) and 4D-consistent human images from sparse-view inputs. We propose a novel sliding iterative denoising strategy to enhance spatial and temporal consistency while maintaining high computational efficiency. To further improve motion accuracy and visual quality, we introduce a 4D pose conditioning mechanism that leverages human skeletons. Our method demonstrates superior performance in capturing fine details and complex human motions from sparse inputs compared to existing state-of-the-art approaches.

Despite promising results, our method still faces certain limitations. First, higher-resolution (4K) videos are not supported due to constraints in base models. Second, our method may struggle with scenarios involving complex human-object interactions. Finally, our current method cannot achieve novel-pose rendering because it requires input videos to constrain the spatial consistency of the generated videos. Addressing these challenges represents interesting directions for future research.

# References

[1] Easymocap - make human motion capture easier. Github, 2021. 6

[2] Kairat Aitpayev and Jaafar Gaber. Creation of 3d human avatar using kinect. *Asian Transactions on Fundamentals of Electronics, Communication & Multimedia*, 1(5):12–24, 2012. 2

[3] Sherwin Bahmani, Xian Liu, Wang Yifan, Ivan Skorokhodov, Victor Rong, Ziwei Liu, Xihui Liu, Jeong Joon Park, Sergey Tulyakov, Gordon Wetzstein, et al. Tc4d: Trajectory-conditioned text-to-4d generation. In *European Conference on Computer Vision*, pages 53–72. Springer, 2024. 3

[4] Sherwin Bahmani, Ivan Skorokhodov, Victor Rong, Gordon Wetzstein, Leonidas Guibas, Peter Wonka, Sergey Tulyakov, Jeong Joon Park, Andrea Tagliasacchi, and David B Lindell. 4d-fy: Text-to-4d generation using hybrid score distillation sampling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7996–8006, 2024. 3

[5] Federica Bogo, Michael J Black, Matthew Loper, and Javier Romero. Detailed full-body reconstructions of moving people from monocular rgb-d sequences. In *Proceedings of the IEEE international conference on computer vision*, pages 2300–2308, 2015. 2

[6] Tim Brooks, Bill Peebles, Connor Holmes, Will DePue, Yufei Guo, Li Jing, David Schnurr, Joe Taylor, Troy Luhman, Eric Luhman, et al. Video generation models as world simulators. 2024. *URL https://openai. com/research/video-generation-models-as-world-simulators*, 3:1, 2024. 3

[7] Ang Cao and Justin Johnson. Hexplane: A fast representation for dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 130–141, 2023. 2

[8] David Charatan, Sizhe Li, Andrea Tagliasacchi, and Vincent Sitzmann. pixelsplat: 3d gaussian splats from image pairs for scalable generalizable 3d reconstruction. In *CVPR*, 2024. 3

[9] Yuedong Chen, Haofei Xu, Chuanxia Zheng, Bohan Zhuang, Marc Pollefeys, Andreas Geiger, Tat-Jen Cham, and Jianfei Cai. Mvsplat: Efficient 3d gaussian splatting from sparse multi-view images. *arXiv preprint arXiv:2403.14627*, 2024. 3

[10] Wei Cheng, Ruixiang Chen, Wanqi Yin, Siming Fan, Keyu Chen, Honglin He, Huiwen Luo, Zhongang Cai, Jingbo Wang, Yang Gao, Zhengming Yu, Zhengyu Lin, Daxuan Ren, Lei Yang, Ziwei Liu, Chen Change Loy, Chen Qian, Wayne Wu, Dahua Lin, Bo Dai, and Kwan-Yee Lin. Dna-rendering: A diverse neural actor repository for high-fidelity human-centric rendering. *arXiv preprint*, arXiv:2307.10173, 2023. 2, 6, 7, 8, 1

[11] Alvaro Collet, Ming Chuang, Pat Sweeney, Don Gillett, Dennis Evseev, David Calabrese, Hugues Hoppe, Adam Kirk, and Steve Sullivan. High-quality streamable free-viewpoint video. *ACM Transactions on Graphics (ToG)*, 34(4):1–13, 2015. 2

[12] Matt Deitke, Dustin Schwenk, Jordi Salvador, Luca Weihs, Oscar Michel, Eli VanderBilt, Ludwig Schmidt, Kiana Ehsani, Aniruddha Kembhavi, and Ali Farhadi. Objaverse: A universe of annotated 3d objects. *arXiv preprint arXiv:2212.08051*, 2022. 2

[13] Matt Deitke, Ruoshi Liu, Matthew Wallingford, Huong Ngo, Oscar Michel, Aditya Kusupati, Alan Fan, Christian Laforte, Vikram Voleti, Samir Yitzhak Gadre, Eli VanderBilt, Aniruddha Kembhavi, Carl Vondrick, Georgia Gkioxari, Kiana Ehsani, Ludwig Schmidt, and Ali Farhadi. Objaverse-xl: A universe of 10m+ 3d objects. *arXiv preprint arXiv:2307.05663*, 2023. 2

[14] Yuanxing Duan, Fangyin Wei, Qiyu Dai, Yuhang He, Wenzheng Chen, and Baoquan Chen. 4d-rotor gaussian splatting: towards efficient novel view synthesis for dynamic scenes. In *ACM SIGGRAPH 2024 Conference Papers*, pages 1–11, 2024. 2

[15] Jiemin Fang, Taoran Yi, Xinggang Wang, Lingxi Xie, Xiaopeng Zhang, Wenyu Liu, Matthias Nießner, and Qi Tian. Fast dynamic radiance fields with time-aware neural voxels. In *SIGGRAPH Asia 2022 Conference Papers*, pages 1–9, 2022. 3

[16] Sara Fridovich-Keil, Giacomo Meanti, Frederik Rahbæk Warburg, Benjamin Recht, and Angjoo Kanazawa. K-planes: Explicit radiance fields in space, time, and appearance. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12479–12488, 2023. 2

[17] Silvano Galliani, Katrin Lasinger, and Konrad Schindler. Massively parallel multiview stereopsis by surface normal diffusion. 2015. 2

[18] Ruiqi Gao, Aleksander Holynski, Philipp Henzler, Arthur Brussee, Ricardo Martin-Brualla, Pratul P. Srinivasan, Jonathan T. Barron, and Ben Poole. Cat3d: Create anything in 3d with multi-view diffusion models. In *NeurIPS 2024*, 2024. 3, 1

[19] Oliver Grau. Studio production system for dynamic 3d content. In *Visual Communications and Image Processing 2003*, pages 80–89. SPIE, 2003. 2

[20] Lukas Höllein, Aljaž Božič, Norman Müller, David Novotny, Hung-Yu Tseng, Christian Richardt, Michael Zollhöfer, and Matthias Nießner. Viewdiff: 3d-consistent image generation with text-to-image models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5043–5052, 2024. 3

[21] Wenyi Hong, Ming Ding, Wendi Zheng, Xinghan Liu, and Jie Tang. Cogvideo: Large-scale pretraining for text-to-video generation via transformers. *arXiv preprint arXiv:2205.15868*, 2022. 3

[22] Shoukang Hu and Ziwei Liu. Gauhuman: Articulated gaussian splatting from monocular human videos. *arXiv preprint arXiv:*, 2023. 3, 6, 7

[23] Mustafa Işık, Martin Rünz, Markos Georgopoulos, Taras Khakhulin, Jonathan Starck, Lourdes Agapito, and Matthias Nießner. Humanrf: High-fidelity neural radiance fields for humans in motion. *ACM Transactions on Graphics (TOG)*, 42(4):1–12, 2023. 6, 7, 8, 2

[24] Yan-Bin Jia. Plücker coordinates for lines in the space. *Problem Solver Techniques for Applied Computer Science, Com-S-477/577 Course Handout*, 3, 2020. 3

[25] Takeo Kanade, Peter Rander, and PJ Narayanan. Virtualized reality: Constructing virtual worlds from real scenes. *IEEE multimedia*, 4(1):34–47, 1997. 2

[26] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics*, 42 (4), 2023. 2, 1

[27] Rawal Khirodkar, Timur Bagautdinov, Julieta Martinez, Su Zhaoen, Austin James, Peter Selednik, Stuart Anderson, and Shunsuke Saito. Sapiens: Foundation for human vision models. *arXiv preprint arXiv:2408.12569*, 2024. 5, 1

[28] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 1

[29] Weijie Kong, Qi Tian, Zijian Zhang, Rox Min, Zuozhuo Dai, Jin Zhou, Jiangfeng Xiong, Xin Li, Bo Wu, Jianwei Zhang, et al. Hunyuanvideo: A systematic framework for large video generative models. *arXiv preprint arXiv:2412.03603*, 2024. 3

[30] Kiriakos N Kutulakos and Steven M Seitz. A theory of shape by space carving. *International journal of computer vision*, 38:199–218, 2000. 1

[31] Bing Li, Cheng Zheng, Wenxuan Zhu, Jinjie Mai, Biao Zhang, Peter Wonka, and Bernard Ghanem. Vivid-zoo: Multi-view video generation with diffusion model. *Advances in Neural Information Processing Systems*, 37:62189–62222, 2025. 3

[32] Tianye Li, Mira Slavcheva, Michael Zollhoefer, Simon Green, Christoph Lassner, Changil Kim, Tanner Schmidt, Steven Lovegrove, Michael Goesele, Richard Newcombe, et al. Neural 3d video synthesis from multi-view video. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5521–5531, 2022. 2

[33] Zhe Li, Zerong Zheng, Lizhen Wang, and Yebin Liu. Animatable gaussians: Learning pose-dependent gaussian maps for high-fidelity human avatar modeling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024. 3

[34] Haotong Lin, Sida Peng, Zhen Xu, Yunzhi Yan, Qing Shuai, Hujun Bao, and Xiaowei Zhou. Efficient neural radiance fields for interactive free-viewpoint video. In *SIGGRAPH Asia Conference Proceedings*, 2022. 3

[35] Shanchuan Lin, Andrey Ryabtsev, Soumyadip Sengupta, Brian L Curless, Steven M Seitz, and Ira Kemelmacher-Shlizerman. Real-time high-resolution background matting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8762–8771, 2021. 1

[36] Lahav Lipson, Zachary Teed, and Jia Deng. Raft-stereo: Multilevel recurrent field transforms for stereo matching. In *2021 International Conference on 3D Vision (3DV)*, pages 218–227. IEEE, 2021. 3

[37] Ruoshi Liu, Rundi Wu, Basile Van Hoorick, Pavel Tokmakov, Sergey Zakharov, and Carl Vondrick. Zero-1-to-3: Zero-shot one image to 3d object, 2023. 3

[38] Tianqi Liu, Guangcong Wang, Shoukang Hu, Liao Shen, Xinyi Ye, Yuhang Zang, Zhiguo Cao, Wei Li, and Ziwei Liu. Mvsgaussian: Fast generalizable gaussian splatting reconstruction from multi-view stereo. In *European Conference on Computer Vision*, pages 37–53. Springer, 2024. 3

[39] Yuan Liu, Cheng Lin, Zijiao Zeng, Xiaoxiao Long, Lingjie Liu, Taku Komura, and Wenping Wang. Syncdreamer: Generating multiview-consistent images from a single-view image. *arXiv preprint arXiv:2309.03453*, 2023. 3

[40] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J Black. Smpl: A skinned multi-person linear model. In *Seminal Graphics Papers: Pushing the Boundaries, Volume 2*, pages 851–866. 2023. 2

[41] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver: A fast ode solver for diffusion probabilistic model sampling in around 10 steps. *arXiv preprint arXiv:2206.00927*, 2022. 1

[42] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM*, 65(1):99–106, 2021. 2

[43] Richard A Newcombe, Dieter Fox, and Steven M Seitz. Dynamicfusion: Reconstruction and tracking of non-rigid scenes in real-time. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 343–352, 2015. 2

[44] Zijie Pan, Zeyu Yang, Xiatian Zhu, and Li Zhang. Efficient4d: Fast dynamic 3d object generation from a single-view video. *arXiv preprint arXiv:2401.08742*, 2024. 3

[45] Keunhong Park, Utkarsh Sinha, Jonathan T. Barron, Sofien Bouaziz, Dan B Goldman, Steven M. Seitz, and Ricardo Martin-Brualla. Nerfies: Deformable neural radiance fields. *ICCV*, 2021. 3

[46] Keunhong Park, Utkarsh Sinha, Peter Hedman, Jonathan T. Barron, Sofien Bouaziz, Dan B Goldman, Ricardo Martin-Brualla, and Steven M. Seitz. Hypernerf: A higher-dimensional representation for topologically varying neural radiance fields. *ACM Trans. Graph.*, 40(6), 2021. 3

[47] Sida Peng, Junting Dong, Qianqian Wang, Shangzhan Zhang, Qing Shuai, Xiaowei Zhou, and Hujun Bao. Animatable neural radiance fields for modeling dynamic human bodies. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14314–14323, 2021. 3

[48] Sida Peng, Yuanqing Zhang, Yinghao Xu, Qianqian Wang, Qing Shuai, Hujun Bao, and Xiaowei Zhou. Neural body: Implicit neural representations with structured latent codes for novel view synthesis of dynamic humans. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9054–9063, 2021. 2

[49] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. *arXiv preprint arXiv:2209.14988*, 2022. 3

[50] Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. D-NeRF: Neural Radiance Fields for Dynamic Scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020. 3

[51] Jiawei Ren, Liang Pan, Jiaxiang Tang, Chi Zhang, Ang Cao, Gang Zeng, and Ziwei Liu. Dreamgaussian4d: Genera-

tive 4d gaussian splatting. *arXiv preprint arXiv:2312.17142*, 2023. 3

[52] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10684–10695, 2022. 1

[53] Johannes Lutz Schönberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 2, 1

[54] Johannes Lutz Schönberger, Enliang Zheng, Marc Pollefeys, and Jan-Michael Frahm. Pixelwise view selection for unstructured multi-view stereo. In *European Conference on Computer Vision (ECCV)*, 2016. 2, 1

[55] Ruizhi Shao, Zerong Zheng, Hanzhang Tu, Boning Liu, Hongwen Zhang, and Yebin Liu. Tensor4d: Efficient neural 4d decomposition for high-fidelity dynamic reconstruction and rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16632–16642, 2023. 2

[56] Ari Shapiro, Andrew Feng, Ruizhe Wang, Hao Li, Mark Bolas, Gerard Medioni, and Evan Suma. Rapid avatar capture and simulation using commodity depth sensors. *Computer Animation and Virtual Worlds*, 25(3-4):201–211, 2014. 2

[57] Yichun Shi, Peng Wang, Jianglong Ye, Mai Long, Kejie Li, and Xiao Yang. Mvdream: Multi-view diffusion for 3d generation. *arXiv preprint arXiv:2308.16512*, 2023. 3

[58] Uriel Singer, Shelly Sheynin, Adam Polyak, Oron Ashual, Iurii Makarov, Filippos Kokkinos, Naman Goyal, Andrea Vedaldi, Devi Parikh, Justin Johnson, et al. Text-to-4d dynamic scene generation. *arXiv preprint arXiv:2301.11280*, 2023. 3

[59] Jonathan Starck and Adrian Hilton. Virtual view synthesis of people from multiple view video sequences. *Graphical Models*, 67(6):600–620, 2005. 2

[60] Jonathan Starck and Adrian Hilton. Surface capture for performance-based animation. *IEEE computer graphics and applications*, 27(3):21–31, 2007. 2

[61] Jing Tong, Jin Zhou, Ligang Liu, Zhigeng Pan, and Hao Yan. Scanning 3d full human bodies using kinects. *IEEE transactions on visualization and computer graphics*, 18(4):643–650, 2012. 2

[62] Patrick von Platen, Suraj Patil, Anton Lozhkov, Pedro Cuenca, Nathan Lambert, Kashif Rasul, Mishig Davaadorj, Dhruv Nair, Sayak Paul, William Berman, Yiyi Xu, Steven Liu, and Thomas Wolf. Diffusers: State-of-the-art diffusion models. https://github.com/huggingface/diffusers, 2022. 1

[63] Liao Wang, Jiakai Zhang, Xinhang Liu, Fuqiang Zhao, Yanshun Zhang, Yingliang Zhang, Minye Wu, Jingyi Yu, and Lan Xu. Fourier plenoctrees for dynamic radiance field rendering in real-time. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13524–13534, 2022. 2

[64] Peng Wang and Yichun Shi. Imagedream: Image-prompt multi-view diffusion for 3d generation. *arXiv preprint arXiv:2312.02201*, 2023. 3

[65] Daniel Watson, Saurabh Saxena, Lala Li, Andrea Tagliasacchi, and David J. Fleet. Controlling space and time with diffusion models. *arXiv preprint arXiv:2407.07860*, 2024. 2, 4

[66] Chung-Yi Weng, Brian Curless, Pratul P Srinivasan, Jonathan T Barron, and Ira Kemelmacher-Shlizerman. Humannerf: Free-viewpoint rendering of moving people from monocular video. In *Proceedings of the IEEE/CVF conference on computer vision and pattern Recognition*, pages 16210–16220, 2022. 2

[67] Guanjun Wu, Taoran Yi, Jiemin Fang, Lingxi Xie, Xiaopeng Zhang, Wei Wei, Wenyu Liu, Qi Tian, and Xinggang Wang. 4d gaussian splatting for real-time dynamic scene rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 20310–20320, 2024. 2

[68] Rundi Wu, Ruiqi Gao, Ben Poole, Alex Trevithick, Changxi Zheng, Jonathan T. Barron, and Aleksander Holynski. Cat4d: Create anything in 4d with multi-view video diffusion models. *arXiv:2411.18613*, 2024. 2, 3, 4, 6, 7

[69] Yiming Xie, Chun-Han Yao, Vikram Voleti, Huaizu Jiang, and Varun Jampani. Sv4d: Dynamic 3d content generation with multi-frame and multi-view consistency. *arXiv preprint arXiv:2407.17470*, 2024. 2, 3

[70] Zhen Xu, Tao Xie, Sida Peng, Haotong Lin, Qing Shuai, Zhiyuan Yu, Guangzhao He, Jiaming Sun, Hujun Bao, and Xiaowei Zhou. Easyvolcap: Accelerating neural volumetric video research. 2023. 5

[71] Zhen Xu, Sida Peng, Chen Geng, Linzhan Mou, Zihan Yan, Jiaming Sun, Hujun Bao, and Xiaowei Zhou. Relightable and animatable neural avatar from sparse-view video. In *CVPR*, 2024. 3

[72] Zhen Xu, Sida Peng, Haotong Lin, Guangzhao He, Jiaming Sun, Yujun Shen, Hujun Bao, and Xiaowei Zhou. 4k4d: Real-time 4d view synthesis at 4k resolution. In *CVPR*, 2024. 2, 1

[73] Zhen Xu, Yinghao Xu, Zhiyuan Yu, Sida Peng, Jiaming Sun, Hujun Bao, and Xiaowei Zhou. Representing long volumetric video with temporal gaussian hierarchy. *ACM Transactions on Graphics*, 43(6), 2024. 2, 4, 5, 6, 7, 1, 3

[74] Zeyu Yang, Zijie Pan, Chun Gu, and Li Zhang. Diffusion2: Dynamic 3d content generation via score composition of orthogonal diffusion models. *arXiv e-prints*, pages arXiv–2404, 2024. 3

[75] Zhuoyi Yang, Jiayan Teng, Wendi Zheng, Ming Ding, Shiyu Huang, Jiazheng Xu, Yuanming Yang, Wenyi Hong, Xiaohan Zhang, Guanyu Feng, et al. Cogvideox: Text-to-video diffusion models with an expert transformer. *arXiv preprint arXiv:2408.06072*, 2024. 3

[76] Zeyu Yang, Hongye Yang, Zijie Pan, and Li Zhang. Real-time photorealistic dynamic scene representation and rendering with 4d gaussian splatting. In *International Conference on Learning Representations (ICLR)*, 2024. 2, 1

[77] Yao Yao, Zixin Luo, Shiwei Li, Tian Fang, and Long Quan. Mvsnet: Depth inference for unstructured multi-view stereo. In *Proceedings of the European conference on computer vision (ECCV)*, pages 767–783, 2018. 3

[78] Yao Yao, Zixin Luo, Shiwei Li, Tianwei Shen, Tian Fang, and Long Quan. Recurrent mvsnet for high-resolution multi-view stereo depth inference. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5525–5534, 2019. 3

[79] Yuyang Yin, Dejia Xu, Zhangyang Wang, Yao Zhao, and Yunchao Wei. 4dgen: Grounded 4d content generation with spatial-temporal consistency. *arXiv preprint arXiv:2312.17225*, 2023. 3

[80] Yifei Zeng, Yanqin Jiang, Siyu Zhu, Yuanxun Lu, Youtian Lin, Hao Zhu, Weiming Hu, Xun Cao, and Yao Yao. Stag4d: Spatial-temporal anchored generative 4d gaussians. In *European Conference on Computer Vision*, pages 163–179. Springer, 2024. 3

[81] Haiyu Zhang, Xinyuan Chen, Yaohui Wang, Xihui Liu, Yunhong Wang, and Yu Qiao. 4diffusion: Multi-view video diffusion model for 4d generation. *Advances in Neural Information Processing Systems*, 37:15272–15295, 2025. 3

[82] Jason Y Zhang, Amy Lin, Moneish Kumar, Tzu-Hsuan Yang, Deva Ramanan, and Shubham Tulsiani. Cameras as rays: Pose estimation via ray diffusion. *arXiv preprint arXiv:2402.14817*, 2024. 4

[83] Yuyang Zhao, Zhiwen Yan, Enze Xie, Lanqing Hong, Zhenguo Li, and Gim Hee Lee. Animate124: Animating one image to 4d dynamic scene. *arXiv preprint arXiv:2311.14603*, 2023. 3

[84] Peng Zheng, Dehong Gao, Deng-Ping Fan, Li Liu, Jorma Laaksonen, Wanli Ouyang, and Nicu Sebe. Bilateral reference for high-resolution dichotomous image segmentation. *CAAI Artificial Intelligence Research*, 3:9150038, 2024. 1

[85] Shunyuan Zheng, Boyao Zhou, Ruizhi Shao, Boning Liu, Shengping Zhang, Liqiang Nie, and Yebin Liu. Gps-gaussian: Generalizable pixel-wise 3d gaussian splatting for real-time human novel view synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024. 3, 6, 7

[86] Yufeng Zheng, Xueting Li, Koki Nagano, Sifei Liu, Otmar Hilliges, and Shalini De Mello. A unified approach for text- and image-guided 4d scene generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7300–7309, 2024. 3

# Diffuman4D: 4D Consistent Human View Synthesis from Sparse-View Videos with Spatio-Temporal Diffusion Models

## Supplementary Material

## A. Model Details

**Training.** We initialize our model from Stable Diffusion 2.1 [52] developed by Diffusers [62]. We apply full-parameter fine-tuning to the diffusion model for 200k iterations with a batch size of 32 and a learning rate of $10^{-5}$ using 32 NVIDIA H20 GPUs. To accommodate input images with the conditions, we expand the input channels of the model's first convolutional layer from 4 to 15, consisting of 4 channels for image latents, 4 for skeleton latents, 6 for Plücker embeddings, and 1 for a conditional mask. Following the previous work [18], the conditional mask is a binary indicator specifying whether an image serves as a conditioning input or a target.

**Sampling.** Following Stable Diffusion 2.1 [52, 62], we use DPM-Solver++ [41] with 24 sampling steps and a classifier-free guidance scale of 3.0. Our sliding iterative denoising strategy takes approximately 2 minutes to generate a sample sequence of length 48 when executed on a single A100 GPU. To improve efficiency, we parallelize the denoising process across 8 A100 GPUs.

**4D reconstrcution.** We employ LongVolcap [73] to reconstruct the 4D human performances from the generated multi-view videos. LongVolcap is an enhanced version of 4DGS [76] with the ability of effectively reconstructing long volumetric videos with a temporal Gaussian hierarchy representation. We initialize the 4D Gaussian primitives with the coarse geometry obtained using the predicted foreground masks and the space carving algorithm [30, 72]. We then follow the same training and evaluation settings as in the original paper [73] to reconstruct the 4D human performances. Specifically, we optimize the model with the Adam optimizer [28] with a learning rate of $1.6e^{-4}$, each model is trained for 100k iterations for a sequence of 7200 frames, which takes around 1 hour on a single NVIDIA RTX 4090 GPU.

## B. Datasets Details

We conduct extensive processing on the original DNA-Rendering [10] dataset to generate high-quality multi-view videos along with additional masks and skeletons for training and evaluation. The processing pipeline includes camera re-calibration, color correction matrices (CCMs) optimization, foreground mask prediction, and human skeleton estimation. We provide detailed descriptions of each step below.

**Camera calibration.** We empirically found that the camera parameters provided in the DNA-Rendering dataset are not accurate enough for reconstruction verified with 3D Gaussian Splatting (3DGS) [26]. In order to achieve pixel-level accuracy, we first re-calibrated the camera parameters using Colmap [53, 54]. We then optimized the color correction matrix for each camera to ensure consistent color across different views.

**Foreground mask prediction.** There are only a few (around 1/6) sequences in the DNA-Rendering dataset that provide ground truth foreground masks. To obtain accurate foreground masks, we leverage three state-of-the-art background removal methods, namely RMBG-2.0 [84], BiRefNet-Portrait [84], and BackgroundMattingV2 [35], and combine their predictions using a voting mechanism to fully leverage the strengths of each approach. Specifically, we found that RMBG-2.0 may incorrectly recognize background objects as foreground, BiRefNet-Portrait may segment small objects as background, and BackgroundMattingV2 may produce inaccurate results for certain human poses. We demonstrate the effectiveness of the voting strategy in Fig. 7.
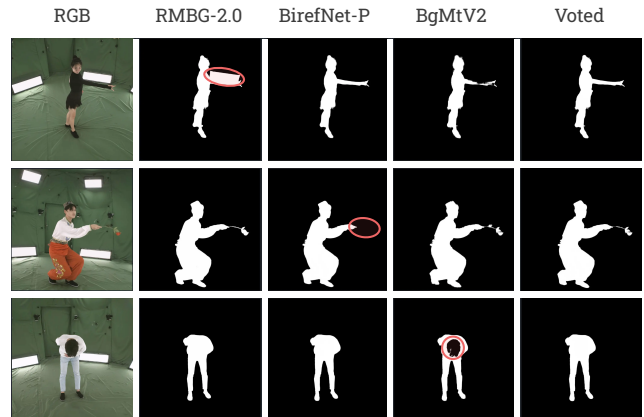


Figure 7. Our voting strategy effectively leverages the strengths of different background removal methods to produce robust foreground masks.

**Human skeleton estimation.** Similar to the foreground mask, only a few sequences have ground truth human skeletons. We thus adopt the state-of-the-art human skeleton estimation model, Sapiens [27], to predict the 2D human skeleton for each frame. We additionally adjust the transparency of the skeleton colors based on the confidence scores of the
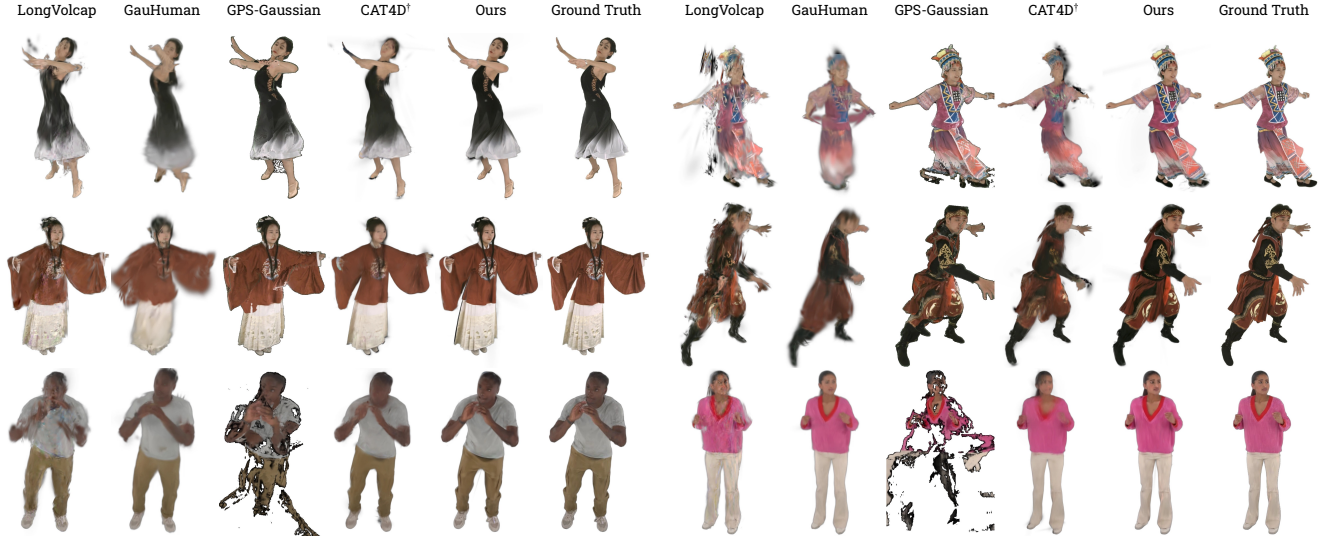
Figure 8. **More qualitative comparisons on DNA-Rendering [10] and ActorsHQ [23].** GPS-Gaussian uses 8 input views while all others use 4 input views, and CAT4D† is our reproduced version. Our Diffuman4D consistently outperforms baselines with higher visual quality and better spatio-temporal consistency.

skeleton, which helps to visualize and encode the uncertainty of the skeleton estimation. After obtaining the 2D human skeletons, we then triangulate them to obtain the 3D human skeleton sequence, which can be further used for projection and evaluation.

We demonstrate the processed data samples in Fig. 9. We plan to release the additionally processed data under the DNA-Rendering open-source license to facilitate future research within the community.

**Dataset filtering.** DNA-Rendering [10] contains many scenes involving human-object interactions, such as writing on a desk, playing guitar, or organizing items. Since the diversity of objects is significantly greater than that of humans, training generative models typically requires extensive object datasets (e.g., Objaverse [12, 13]). To address the relatively limited scale of the DNA-Rendering dataset, we employed the Llama Vision 3.2 model to classify all scenes and filtered out those containing large objects to avoid potential model collapse during training.

Nevertheless, we observe that even though the training dataset does not include objects, our model successfully generalizes to scenes featuring simple objects, such as the basketball player shown in Fig. 10.

## C. Additional Comparisons

**More qualitative results.** We provide additional comparisons with baselines in Fig. 8. Results show that our method consistently outperforms the baselines in terms of visual quality and fine details.

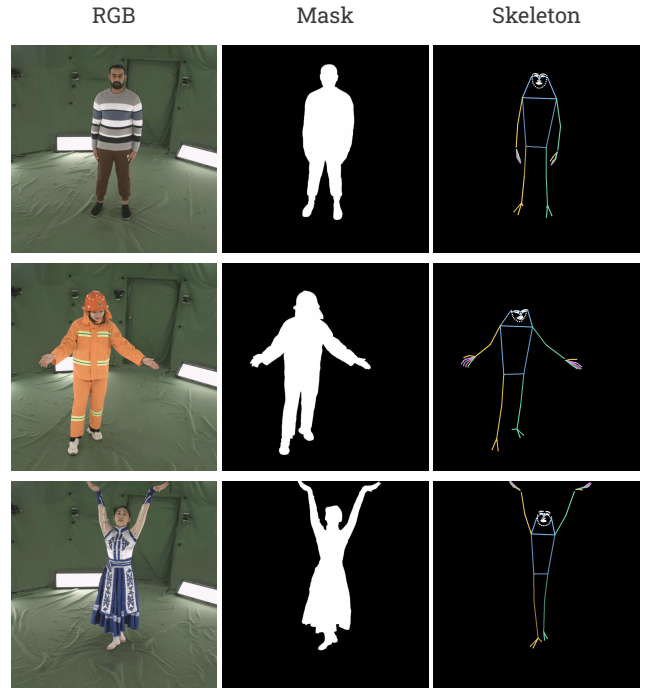**Diffusion generation vs. 4DGS rendering.** Although our



Figure 9. High-quality foreground masks and human skeletons predicted using state-of-the-art methods.

model already supports novel-view synthesis, we choose to optimize a 4DGS model using LongVolcap [73] to enable real-time rendering. As shown in Fig. 10, our model can generate high-fidelity human videos, but they still inevitably exhibit spatio-temporal inconsistencies. Recon-

Figure 10. Qualitative comparisons between novel views generated by our model and those rendered from the 4DGS model reconstructed using LongVolcap [73].

structing a 4DGS model further alleviates these inconsistencies, though at the cost of reduced sharpness compared to the originally generated images.